

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 09-231053

(43)Date of publication of application : 05.09.1997

(51)Int.Cl.

G06F 7/24

G06F 7/36

G06F 12/00

(21)Application number : 08-039799

(71)Applicant : NEC SOFTWARE LTD

(22)Date of filing : 27.02.1996

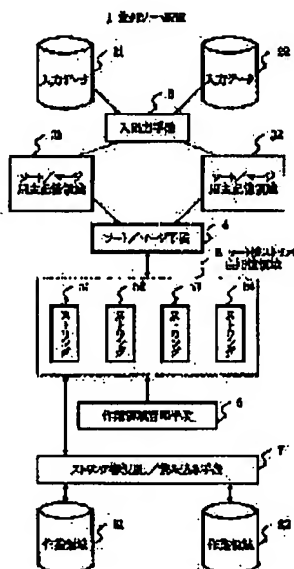
(72)Inventor : KITAZAWA ATSUSHI

## (54) PARALLEL SORT DEVICE

## (57)Abstract:

**PROBLEM TO BE SOLVED:** To shorten the parallel processing time of a parallel sort device by equalizing the numbers of sorted strings which are read out of a secondary storage in every parallel processing unit when the external sorts whose objects data are not completely stored in a main storage are carried out in parallel to each other in a multiprocess.

**SOLUTION:** The sort processing operations are carried out in the sort/merge main storage areas 31 and 32 which are proper to the parallel processes. Then these sort processing results are stored in a shared sorted string main storage area 5. Then the strings are written in the work areas 81 and 82, and a work area management means 6 decides a string replacement block when the area 5 has an overflow. The strings are taken out of the area 5 and sent to the areas 31 and 32 to perform the merging operations after the sort phase is over. Then a string writing/reading means 7 decides the blocks to be read next, and these blocks are read out of the areas 81 and 82.



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平9-231053

(43) 公開日 平成9年(1997)9月5日

| (51) Int.Cl. <sup>8</sup> | 識別記号  | 庁内整理番号 | F I          | 技術表示箇所 |
|---------------------------|-------|--------|--------------|--------|
| G 0 6 F 7/24              |       |        | G 0 6 F 7/24 | M      |
| 7/36                      |       |        | 7/36         |        |
| 12/00                     | 5 1 2 |        | 12/00        | 5 1 2  |

審査請求 有 請求項の数 2 O L (全 7 頁)

(21) 出願番号 特願平8-39799

(22) 出願日 平成8年(1996)2月27日

(71) 出願人 000232092

日本電気ソフトウェア株式会社  
東京都江東区新木場一丁目18番6号

(72) 発明者 北沢 敦

東京都江東区新木場一丁目18番6号 日本  
電気ソフトウェア株式会社内

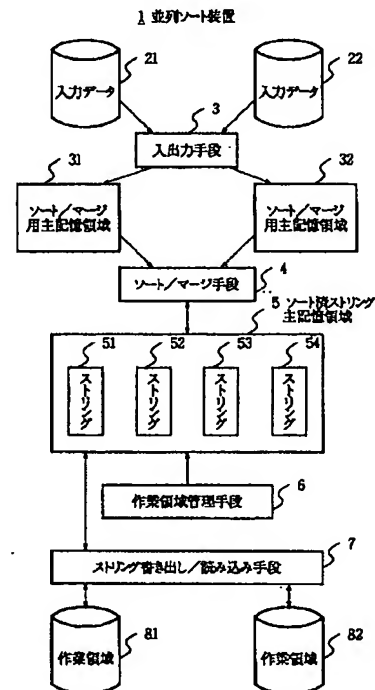
(74) 代理人 弁理士 京本 直樹 (外2名)

(54) 【発明の名称】 並列ソート装置

(57) 【要約】

【課題】 対象データが全て主記憶に納まらない外部ソートをマルチプロセスで並列実行する際に、ソート済みストリングの二次記憶からの読み出し量を各並列処理単位で均等にする事で、並列処理時間を短縮する。

【解決手段】 並列処理プロセスに固有のソートマージ用主記憶領域31、32で、ソート処理を行い、結果を共有のソート済みストリング用主記憶領域5に格納する。作業領域81、82にストリングを書き出し、ソート済みストリング用主記憶領域5がオーバーフローした時、作業領域管理手段6がストリング置き換えブロックを決定し、ソートフェーズ終了後にソート済みストリング用主記憶5からソート/マージ用主記憶31、32にストリングを取り出しマージを行い、ストリング書き出し/読み込み手段7が次に読むべきブロックを決定し、作業領域81、82から読み込みを行う。



## 【特許請求の範囲】

【請求項1】 ソート対象の入力データを格納する第一の二次記憶と、

前記第一の二次記憶から読み出す入力データの所定数ブロックからなるストリングあるいはソート／マージ処理の中間段階で生成されるストリングを保持するソート／マージ用主記憶領域と、

ソート／マージ処理の中間段階で生成されるストリングの格納領域をオーバフローするブロックを退避させる作業領域を有する第二の二次記憶と、

を各プロセスそれぞれで具備し、

前記ソート／マージ用主記憶領域のストリングをソート／マージ処理して生成する新たなストリングを格納する前記格納領域からなるソート済ストリング主記憶領域と、

前記第一の二次記憶から入力データの所定数ブロックからなるストリングを読み出して前記ソート／マージ用主記憶領域に格納する入出力手段と、

前記ソート／マージ用主記憶領域のストリングをソート／マージ処理して生成した新たなストリングを前記ソート済ストリング主記憶領域に格納し、あるいは前記ソート済ストリング主記憶領域のストリングをソート／マージ処理して新たなストリングを前記ソート／マージ用主記憶領域に移送するソート／マージ手段と、を全プロセスで共有してなる並列ソート装置にあって、

前記ソート／マージ手段が前記ソート／マージ用主記憶領域のストリングをソート／マージ処理して生成する新たなストリングのブロックを前記ソート済ストリング主記憶領域のストリングに追加して生じるオーバフローブロックを前記作業領域に置き換えるとき、前記ストリングそれぞれの置き換えられるブロック数を等しくするように制御する作業領域管理手段と、

前記ソート済ストリング主記憶領域のオーバフローブロックを前記作業領域に前記置き換えのため書き出し、あるいは前記ソート済ストリング主記憶領域のストリングをソート／マージ処理するとき、前記置き換えによる空きブロックを、前記作業領域から読み出してストリングを復元するストリング書き出し／読み込み手段と、を備えることを特徴とする並列ソート装置。

【請求項2】 前記プロセスそれぞれをメモリ共有型マルチプロセッサの各プロセッサに割り当てることを特徴とする請求項1記載の並列ソート装置。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】この発明は、ソート／マージ処理の並列実行に関し、特に、ソート対象データが主記憶の共有領域と複数の二次記憶に分離格納されている場合のソーティングに関する。

【0002】関係データベースシステム内部でソーティングが利用される場合、関係データを格納する二次記憶

が分離格納されている場合や、導出された関係データを格納する、中間的に利用される主記憶と二次記憶が分離されている場合がある。

【0003】近年実用化されている並列システムでは、主記憶を処理モジュールが共有するshared everything型のアーキテクチャで、この発明は利用できる。さらに、主記憶を共有しないshared nothing型でも、分散共有メモリアーキテクチャを採用したシステムは、この発明を適用できる。

## 【0004】

【従来の技術】この発明の関連分野の文献を挙げると、

[1] 特開平03-071227

[2] Bitton D., et al: "Parallel Algorithms for the Execution of Relational Database ACM Transactions on Database Systems, Vol 8, No 3, September 1983, pages 324-353.

[3] 樋川英治、渡部栄一: "NonStop SQLのアーキテクチャ", 情報処理, Vol 33, No 12, pages 1436-1440がある。これらの文献を引用して、従来技術を説明する。

【0005】文献[1]によれば、主記憶(容量m)にソート対象となるデータ(容量n)がすべて収まらない様な大容量のソート処理を外部的ソートと呼び、1回のソートフェーズと $\log(n/m)$ 回のマージフェーズをもって構成される[文献1の4頁]。各フェーズでは、その中で完全にソートされた結果をもつ、いくつかのグループが生成する。このグループをストリングと呼ぶこととする。たとえば、ソートフェーズ完了時のストリングの個数は、最大 $n/m$ である。外部ソートの性能は、ソートフェーズとマージフェーズあるいは各々のマージフェーズ間でデータ受け渡しをする作業領域に対するデータ転送時間が全体性能を制限する傾向にある(つまり、作業領域にたいするI/O boundである)

[文献1の8ページ]。[文献1]では、ソートフェーズとマージフェーズの受け渡しを、仮想記憶にマッピングすることで、仮想記憶をサポートするメモリ量が十分にある場合、30%から80%この時間を縮小する方法を開示している。この方法によれば、ソート処理に使用するメモリをソート／マージ処理用と作業領域支援用に等しく分割し、作業領域支援用のメモリは一つのストリングが処理された時点で(すなわちソート／マージ処理領域の全てのデータが処理された時点)、その都度、二次記憶に書き出し、次のマージフェーズでは主記憶に残っているストリングから処理することで作業領域からの読み出し時に発生する二次記憶へのI/O回数を削減するものである。なお二次記憶への書きだしは入力データの読み込みとオーバーラップさせている。

【0006】一方、外部ソートを複数プロセッサの並列マシンに適用した例が、文献[2]に示されている。図5は[文献2の頁334]から引用した並列2ウェイマージソートを説明する図で、SUB OPTIMAL STAGEとOPTIMAL STAGEが一つのデータ源に対する外部ソートに対応する。すなわち外部データ91乃至94をプロセッサ95乃至98でソートし、その結果を最終的に各プロセッサ95乃至98毎に1本のストリング99にマージしている。並列2ウェイマージソートでは、個々の外部ソートの結果をさらに多段階に構成された複数プロセッサでマージすることで(POSTOPTIMAL STAGE)、1本のソートされた結果を得ることができる。

【0007】関係データベースシステムにおける関係代数演算処理を並列処理する場合、文献[3]によれば、関係を複数の二次記憶装置に分離して格納しておき、各々の分離した関係(パーティションと呼ぶ)に対して関係代数演算を施し、その結果をまとめる場合がある。この場合の一般的な考慮として、データの片寄りの問題がある。関係代数演算処理では、ソート処理の入力が通常選択処理等の別の関係代数演算の結果生じるので、データの片寄り、単に各パーティションのデータ量が各々異なるだけでなく、パーティション毎の選択率が異なることによる片寄りもある。この様に並列処理する対象のデータ量に片寄りがある場合、並列ソートのOPTIMAL STAGEとPOSTOPTIMAL STAGE間で最もデータ処理量の多いプロセッサの処理終了を待ち合わせる必要が発生することになり、全体性能を低下させる。

【0008】

【発明が解決しようとする課題】メモリ共有型、あるいはクラスタ型のアーキテクチャを持つ並列マシン上で並列ソート処理を実行する場合の問題は、SUBOPTIMAL STAGE、OPTIMAL STAGEを個々の外部ソートとみなし、単純に文献[1]の方法を適用した場合、各外部ソートで発生する作業領域の二次記憶へのI/O回数が異なり、最も多くのI/Oを必要とする外部ソートの性能で、並列ソート全体の性能が抑制されることである。これはメモリ共有型のアーキテクチャの特徴を生かしていない。他の外部ソートのメモリを融通することで、I/O回数のバランスをとることが可能なはずだからである。

【0009】外部ソートの場合、各フェーズで同一データに関する二次記憶へのI/O回数は高々2度である。1度は入力データの読み出し、あるいは前フェーズの作業領域からの読み出しであり、1度は作業領域へのデータの書き込みである。したがって、I/O回数を削減する方法は、フェーズ数を削減するか、あるいは各フェーズの作業領域へのデータの読み書きを主記憶でバックアップするかのいずれかである。フェーズ数を削減するた

めには、ソート/マージ用に割り当てる主記憶量を調整する必要がある。ソート/マージ用に割り当てられた主記憶はランダムにアクセスする傾向にあり、作業領域用に割り当てられた主記憶は順次アクセスとなる。

【0010】この発明では、並列ソートのSUBOPTIMAL STAGE、OPTIMAL STAGEを外部ソートとみなし、作業領域用に割り当てるメモリを並列ソートの間で融通しあうことで、マージフェーズのストリングで2次記憶から読み出すブロック数を各ストリングで均等にし、並列動作する外部ソートの終了時刻をなるべく同じにする。

【0011】

【課題を解決するための手段】この発明の目的は、関係データベースの間合せ処理の並列化処理で、ソート処理を並列実行する際、二次記憶に対する入出力回数のバラツキで生じる処理時間の片寄りを均等に調整することによって、全体性能を向上させることにある。

【0012】そのため、この発明の、ソート対象の入力データを格納する第一の二次記憶と、前記第一の二次記憶から読み出す入力データの所定数ブロックからなるストリングあるいはソート/マージ処理の中間段階で生成されるストリングを保持するソート/マージ用主記憶領域と、ソート/マージ処理の中間段階で生成されるストリングの格納領域をオーバーフローするブロックを退避させる作業領域を有する第二の二次記憶と、を各プロセスそれぞれで具備し、前記ソート/マージ用主記憶領域のストリングをソート/マージ処理して生成する新たなストリングを格納する前記格納領域からなるソート済ストリング主記憶領域と、前記第一の二次記憶から入力データの所定数ブロックからなるストリングを読み出して前記ソート/マージ用主記憶領域に格納する入出力手段と、前記ソート/マージ用主記憶領域のストリングをソート/マージ処理して生成した新たなストリングを前記ソート済ストリング主記憶領域に格納し、あるいは前記ソート済ストリング主記憶領域のストリングをソート/マージ処理して新たなストリングを前記ソート/マージ用主記憶領域に移送するソート/マージ手段と、を全プロセスで共有してなる並列ソート装置にあって、前記ソート/マージ手段が前記ソート/マージ用主記憶領域のストリングをソート/マージ処理して生成する新たなストリングのブロックを前記ソート済ストリング主記憶領域のストリングに追加して生じるオーバーフローブロックを前記作業領域に置き換えるとき、前記ストリングそれぞれの置き換えられるブロック数を等しくするように制御する作業領域管理手段と、前記ソート済ストリング主記憶領域のオーバーフローブロックを前記作業領域に前記置き換えのため書き出し、あるいは前記ソート済ストリング主記憶領域のストリングをソート/マージ処理するとき、前記置き換えによる空きブロックを、前記作業領域から読み出してストリングを復元するストリング書き

出し／読み込み手段と、を備えることを特徴とする。

【0013】ソートの対象となるデータは、二次記憶の入力データからソート／マージ用主記憶領域に入り切る所定数のブロックをストリングとして取り出され、その中で内部ソートされる。この処理は、別々のプロセス（あるいはプロセッサ）によって並列に実施される。いったん内部ソートが完了して生成される新たなストリングは個々のプロセスでソート済ストリング用主記憶領域に転送される。ソート済ストリング用主記憶のストリングのブロックは、二次記憶の作業領域に転送されて対応付けられる。ソート／マージ用主記憶領域に各々のプロセスで新たな入力データが読み出されるが、この際の待ち時間には、ストリング書きだし／読み出し手段がソート／マージ用主記憶領域のストリング中のブロックを書き出せるだけソート済ストリング主記憶領域と二次記憶の作業領域に書き出す。前記動作中に、転送すべきソート済ストリング用主記憶に空きがなくなった場合、作業領域管理手段が、各ストリングで置き換えブロック数が等しくなる様に置き換えるべきブロックを決定する。この決定方法については、実施例で説明する。

【0014】マージフェーズでは、ソート済ストリング用主記憶領域の順次ソート済のストリングを順に取り出し、マージ用主記憶領域に転送してマージ処理を行う。これは、複数のプロセスで並列に実施される。この動作を各プロセスが作成した全てのストリングに対して実施する。マージ処理が完了したものは、従来技術で説明したPOSTOPTIMAL STAGEの処理に渡される。これは、同一マシン内で実施してもよいし、別マシンに転送して実施してもよい。

【0015】上記動作中に、ソート済ストリング用主記憶のブロックに空きができた場合、ストリング書き出し／読み込み手段は二次記憶の作業領域から空きになったブロックと同じストリングで次に利用するブロックの読み出し要求を出す。どのブロックを読み出すかについての方法を実施例で説明する。二次記憶の作業領域からブロックの読み出しは、マージ処理実行中に行われる。マージ処理は、ソート済ストリング主記憶領域中のストリングがなくなった場合、ストリング書き出し／読み出し手段が要求した読み出し要求の終了を待ち合わせる。この動作によって、各プロセスが二次記憶の作業領域から読み出すブロック数を均等にして、プロセス間で外部ソートに必要な処理時間を合わせることが出来る。

【0016】

【発明の実施の形態】次に、この発明について、図面を参照して説明する。

【0017】この発明の一実施例の構成を示す図1を参照すると、並列ソート装置1が2つのプロセスで動作する構成を例示し、ソートの入力となる入力データ21、22と、入力データをプロセス対応のソートマージ用主記憶領域31、32に読み込むための入出力手段3と、

読み込んだ入力データをソート／マージ用主記憶領域31、32上でソートあるいはマージして新しいストリングを生成するソート／マージ手段4と、ソート／マージ手段4によってソートされた結果であるストリング51乃至54を蓄えるソート済ストリング用主記憶領域5と、入力データの量が多く、ソート済ストリング領域5がオーバフローして、該主記憶領域5が不足した場合に備えて、ストリング51乃至54でどのブロックを新たなブロックの置き換え用として選択するかを決定する作業領域管理手段6と、作業領域管理手段6が決定したブロックをプロセス対応の作業領域81、82に書き出すストリング書き出し／読み込み手段7と、を備える。

【0018】この実施例では、各並列処理のプロセスがソート結果を1本のストリングとしている。これは、従来技術で述べた並列ソートのSUBOPTIMAL STAGE、OPTIMAL STAGEに相当する。最終的に各プロセスで処理したソート済ストリングから1本のソート結果を得るには、さらにその結果をマージする必要がある。入力データ21、22とソート／マージ用主記憶領域31、32と、作業領域81、82は並列処理単位のプロセス毎に存在し、ソート済ストリング用主記憶領域5のみが、並列処理のプロセスに共通した領域であるものとしている。作業領域81、82や、ソート／マージ用主記憶領域31、32を並列処理単位に共通にする変更を加えることは可能である。

【0019】次に、この実施例のソートマージ用主記憶領域31、32と、ソート済ストリング用主記憶領域5のサイズの決定方法を説明する。更に作業領域管理手段6がソート済ストリング用主記憶領域5上のストリング51乃至54のどのブロックを置き換えるかの決定方法を説明し、更にソート済ストリング用主記憶領域5に空きブロックが発生した場合、ストリング書き出し／読み出し手段7がどのブロックを二次記憶の作業領域81、82から復元するかの決定方法を説明する。

【0020】先ず主記憶領域のサイズの決定方法を説明する。ソート／マージ用主記憶はプロセス固有の領域であり、サイズは各プロセスで同じにする。各プロセスが処理する入力データサイズを $n$ とし、ソート用主記憶のサイズを $m$ とすると、一つのプロセスが生成するソート済ストリングの数は最大で $n/m$ となる。一方、マージ処理が必要とするメモリ量は、処理対象のストリング数に依存して増加する。簡単のために、マージ対象のストリングは、そのプロセスが生成したストリングであるとする。マージ処理に必要な主記憶量を $f(n/m)$ で表わすと、

$$m = f(n/m)$$

なるメモリ量 $m$ を決定することが望ましい。たとえば、 $f(n/m) = n/m$ のケースを考えると、上式は $m = n/m$ となり、 $m = \text{root}(n)$ となる。上述のように、ソート／マージ用主記憶領域サイズは、マージアル

ゴリズムで使用するメモリ量と処理の入力データとから決定することができる。入力データのデータ量の半分の桁数のメモリ量を割り当てることになる。利用可能な主記憶のうち、残りの主記憶をソート済ストリング用主記憶領域に割り当てる。

【0021】次に、作業領域管理手段6のブロック置き換えの動作の流れ図3を参照し、ブロック置き換えの動作を説明する図2(a)、図2(b)、図2(c)、図2(d)を援用して、作業領域管理手段6による置き換えアルゴリズムを説明する。図2(a)乃至図2(d)では、ソート処理結果としてストリング51乃至54の4本が生成され、ソート済ストリング用主記憶領域5に渡されるものとしている。ソート済ストリング用主記憶領域5に存在するブロックは全部で12ブロック分であるものとしている。また12ブロックがオーバフローして、ブロックの置き換えが発生し、ソート済ストリング用主記憶領域5から作業領域81、82に置き換えで追い出されたブロックは黒の塗り潰しで示すことにする。図2(c)、図2(d)で、置き換え数とは、この黒塗りのブロック数である。また、残りとは、図2(a)、(b)で網点で示されているソート済ストリング主記憶領域5の12ブロック内に残っている数である。空白のブロックは、今後処理で生成される予定のストリングのブロックである。ストリング51とストリング52はプロセス15がソート処理したストリングであり、後のマージ処理でもプロセス15が使用する。ストリング53とストリング54はプロセス16がソート処理したストリングであり、後のマージ処理ではプロセス16が使用する。

【0022】先ず図2(a)は、ストリング53の生成時に領域がオーバフローし、初めて置き換えが必要になった時点でソート済ストリング主記憶領域5の12ブロック内にストリングが納められた状況を示している。この時点で、ストリング51、52、53の全ての置き換え数は0であるが、生成したブロックを該領域5に格納するとき、図3によれば、自分以外(つまり処理対象となっているストリング53以外)を優先的に選択するので、置き換え対象としてストリング51が選ばれる(ステップ61)。ストリング51の1ブロックを置き換えた場合にソート済ストリング用主記憶に残るブロックが存在する(1以上である)ので、このストリングの最後に転送された部分(図2の513)が置き換えの対象となったことを示している(ステップ62のno)。なお、ブロック1個を残すのは、全てのソートストリングが生成された後のストリングのマージ処理を実行するためには、全てのストリングの最も小さい値を含むブロックがソート済ストリング用主記憶領域5に存在する必要があるからである。図3のアルゴリズムにしたがって、4本全てのストリングの処理が終わった時点の状態を図2(b)に示す。ストリング54のブロックは、自身の

ストリング書き出し中に発生した黒の塗り潰し部分512、513、522、523、536、537、542、546の置き換えによって、12ブロック内の網点511、521、531乃至535、541、543乃至545、547に納められている。この様な置き換えは、図3のステップ62の条件から、図2のブロック541では発生しない。図2(b)によれば最終的に全てのストリングで置き換え数が2となり、読み出すべきブロック数が同じになる様に調整されていることがわかる。このことは、マージ処理時にプロセス15およびプロセス16で必要な作業領域81、82からの読みだし回数が等しくなり、各プロセスでのマージ処理終了時間が均等になることを意味する。ソート済ストリング主記憶領域5にストリングが生成されて、ブロックを置き換えるためには、ソート済ストリング用主記憶5の置き換え対象となるブロックをストリング書き出し/読みだし手段7で作業領域81、82に書き出す必要がある。書き出すタイミングを、置き換えが発生するタイミングとしたが、事前に書き出ししておくこともできる。たとえば、最初に置き換えが発生する一つ前(つまり、ソート済ストリング用主記憶領域に1ブロックの空きがある時点で)に、図3のアルゴリズムを適用し、次に置き換えが発生した場合に置き換えるブロックの書き出し要求を出しておくことで、入力データからの読み出しおよびソート処理と作業領域81、82へのブロックの書き出し処理をオーバーラップさせて、処理時間を短縮することができる。マージフェーズでソート済ストリング領域5の或るブロックの処理が完了し、そのブロックが再利用可能となった場合、ストリング書き出し/読み込み手段7がどのブロックを作業領域81、82から読み出すかを説明する。

【0023】先ず図2(b)がマージ処理の開始時点である。ストリング書き出し/読み込み手段7のブロック選択基準を示すフローチャートである図4によれば図2(b)で、プロセス16のマージ処理が進行し、ブロック531が空になったとする。ここで、図4の1のフローチャートに従うと、プロセス56自身のマージ処理で、空きのブロック531に読み出す作業領域のブロックは、図2(b)の536、542から置き換えられたブロックで、このうち542は距離1であり、536は距離5であることから(ステップ71、72)、ストリング書き出し/読み込み手段7はブロック542を読み出す(ステップ73)。この間、マージ処理は、図2(b)のブロック541、532の間で継続する。一方、プロセス15のマージ処理が進行しブロック511が空になったとすると、次に読み出すべきブロックは図2(b)の512、522で置き換えられたブロックであり、各々距離は1なので、自分のストリングの次のブロックである512が読み出しの対象となる。

【0024】なお、この実施例では、ダブルバッファに

よる先読みを考慮していないが、図3のステップ2で、残すバッファ数を調整することで、ダブルバッファによる先読みを実施することもできる。たとえば、図3のステップ2の判定条件を、「残りが1以下または、そのプロセスの先読み用バッファが1以下」あるいは、「残りが2以下」などと変更すればよい。この場合、図4のアルゴリズムはまったく変わらない。またこの実施例において、プロセスを共有メモリ型マルチプロセッサの各プロセスに割り当てる変更ができることは明らかなである。

#### 【0025】

【発明の効果】以上説明した様に、この発明によれば、並列に実行する外部ソートのマージフェーズで読み出すべきブロック数を均等になるように調整することによって、マージフェーズの終了時刻を各プロセスで同じにして、ソート／マージ処理全体の処理時間を短縮できる。たとえば、図2(b)の例では、プロセス毎に別々の作業領域用の二次記憶を持つとすると、二次記憶へのI/O回数は各プロセスで4回ずつである。これに対して、各プロセスで均等にソート済ストリング用主記憶を割り当てたとすると、プロセス15は全てメモリに収まるために、I/O回数は0であり、プロセス16では8回のI/Oが必要となり、この発明の装置の2倍のI/O回数のプロセスが全体の処理時間を低下させる。

【0026】また、この発明によれば、空きブロックに読み出すブロックを決定する際に、現在使用しているブロックからの距離を利用することで、最も必要度の高いデータを読み出しの対象とすることができ、I/O中にマージ処理を並列して行うことができる。たとえば、図2(b)の例では、ブロック531が空となった場合にブロック536を読み出すと、次にブロック542が必

要になった場合に、ブロック542の読み出し完了までマージ処理を待つ必要がある。これは、ブロック542がブロック536よりも前に必要となる可能性が高いから各ストリングが均等なソートキーの分布を持つ場合には有効になる。

#### 【図面の簡単な説明】

【図1】この発明の一実施例の構成を示す図である。

【図2】分図(a)は、ストリングに新たなブロックをつなぐとき、置き換えるブロックを決定する説明図、分図(b)は、ソート済ストリングをマージするとき、読み出すブロックを決定する説明図、分図(c)は、分図(a)の時点でのストリングの状況を説明する図、分図(d)は分図(b)の時点でのストリングの状況を説明する図である。

【図3】ソート済ストリングを主記憶領域で置き換える際のフローチャートである。

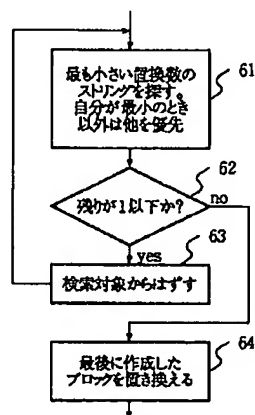
【図4】ソート済ストリング用主記憶領域に空ができた場合の二次記憶の作業領域からの読み出しのフローチャートである。

【図5】従来の並列ソートを示す図である。

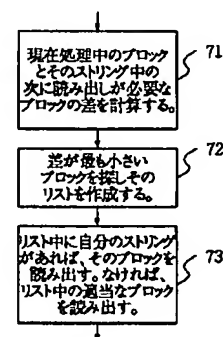
#### 【符号の説明】

- 1 並列ソート装置
- 3 入出力手段
- 4 ソート／マージ手段
- 5 ソート済ストリング主記憶領域
- 6 作業領域管理手段
- 7 ストリング書き出し／読み込み手段
- 21, 22 入力データ
- 31, 32 ソート／マージ用主記憶領域
- 81, 82 作業領域

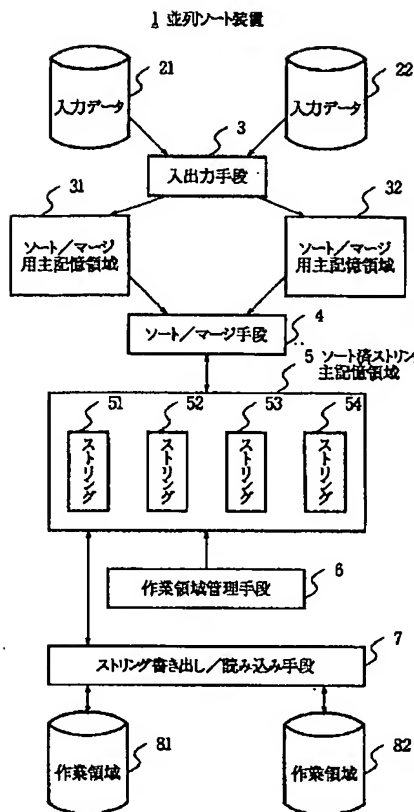
【図3】



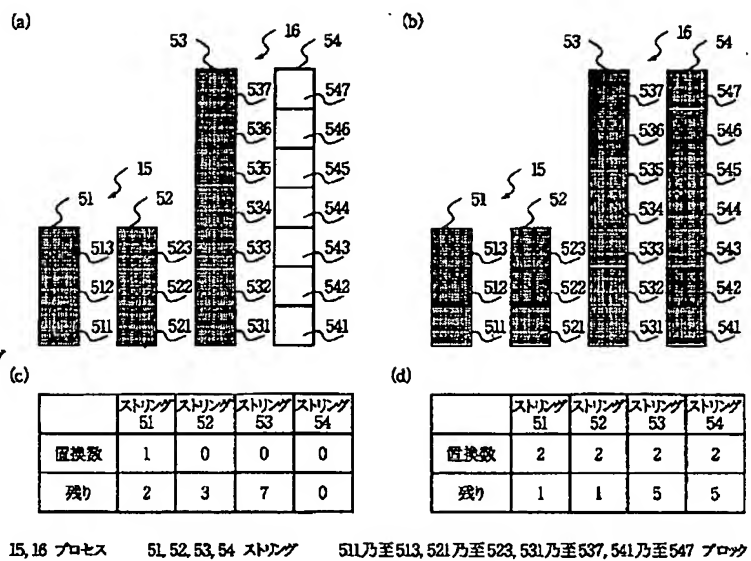
【図4】



【図1】



【図2】



【図5】

